



# Data-Driven Behavioral Analytics: Observations, Representations and Models

Meng Jiang (UIUC)

Peng Cui (Tsinghua)

Jiawei Han (UIUC)

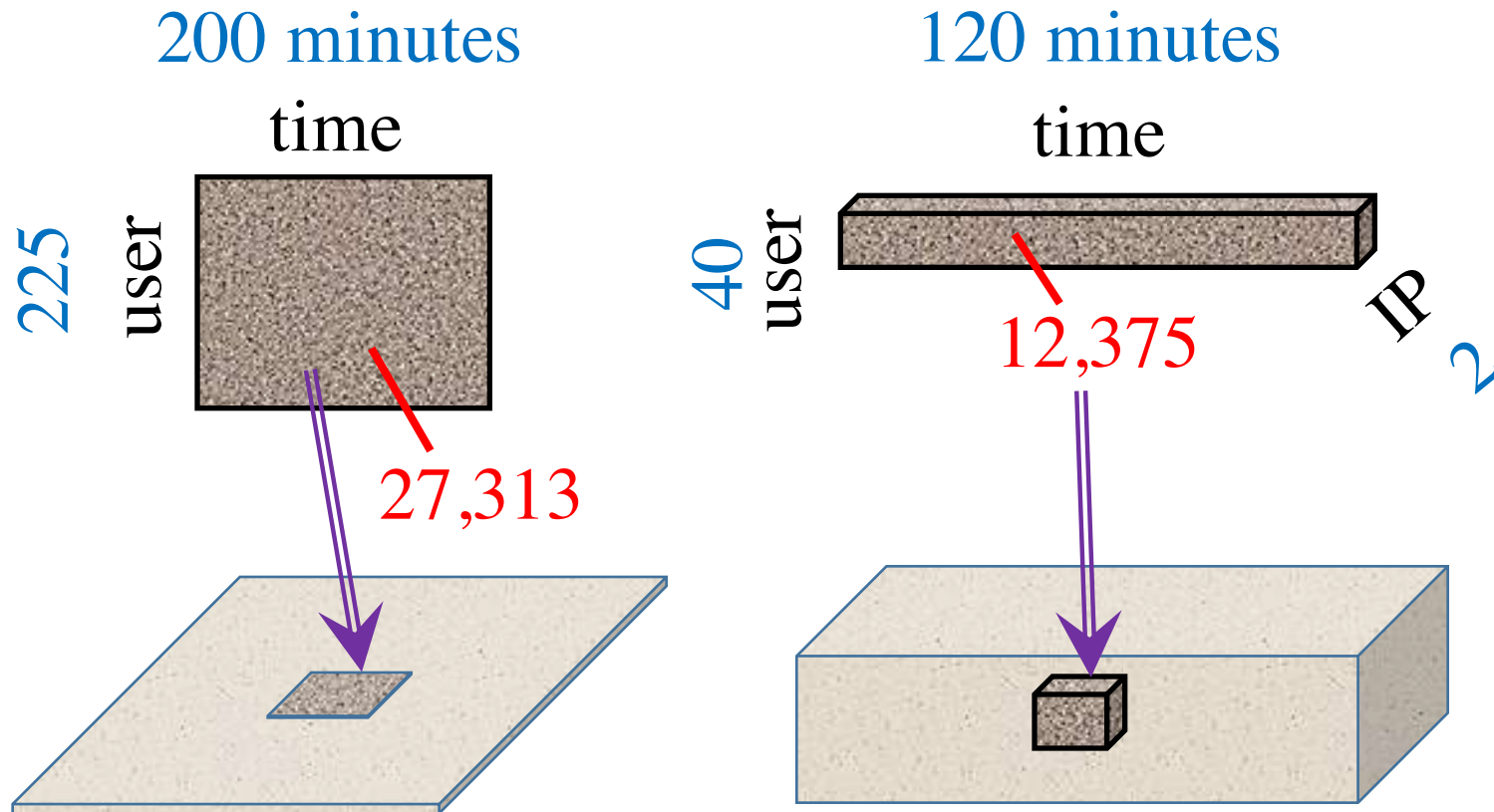
<http://www.meng-jiang.com/tutorial-cikm16.html>

# Observation: Spatiotemporal Contexts

<b>Dataset</b>	<b>Dimension/Mode</b>				<b>Mass</b>
<b>Weibo's Retweeting</b>	<b>User</b>	<b>Root ID</b>	<b>IP</b>	<b>Time (min)</b>	<b>#retweet</b>
	29.5M	19.8M	27.8M	56.9K	211.7M
<b>Weibo's Trending (Hashtag)</b>	<b>User</b>	<b>Hashtag</b>	<b>IP</b>	<b>Time (min)</b>	<b>#tweet</b>
	81.2M	1.6M	47.7M	56.9K	276.9M
<b>Network attacks (LBNL)</b>	<b>Src-IP</b>	<b>Dest-IP</b>	<b>Port</b>	<b>Time (sec)</b>	<b>#packet</b>
	2,345	2,355	6,055	3,610	230,836

Jiang et al. **A General Suspiciousness Metric for Dense Blocks in Multimodal Data.** *ICDM*, 2015.

# Dense Block Indicates Suspiciousness

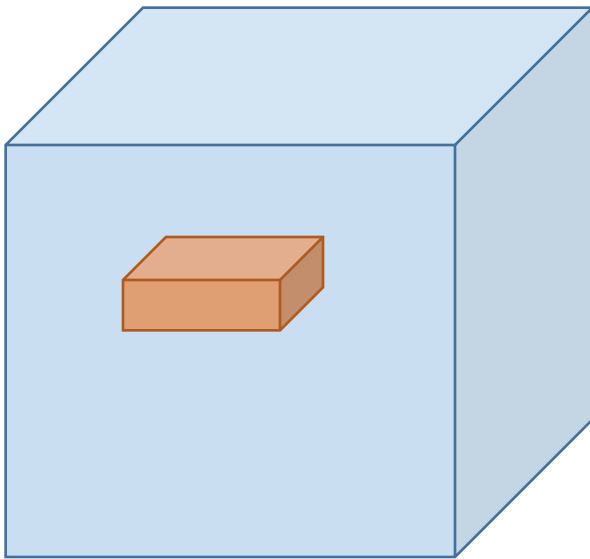


**Q:** Which is more suspicious?

We need a metric to evaluate the suspiciousness.

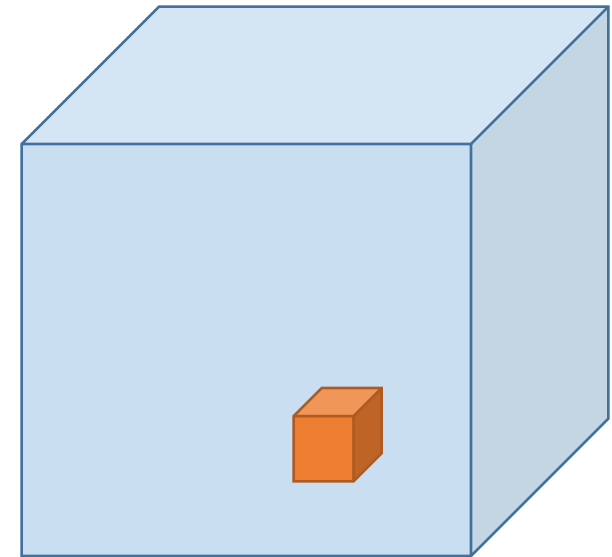
# Criteria for Suspiciousness Metric

What properties are required of a good metric?



$$N_1 \times N_2 \times N_3$$

Count data with  
total “mass”  $C$



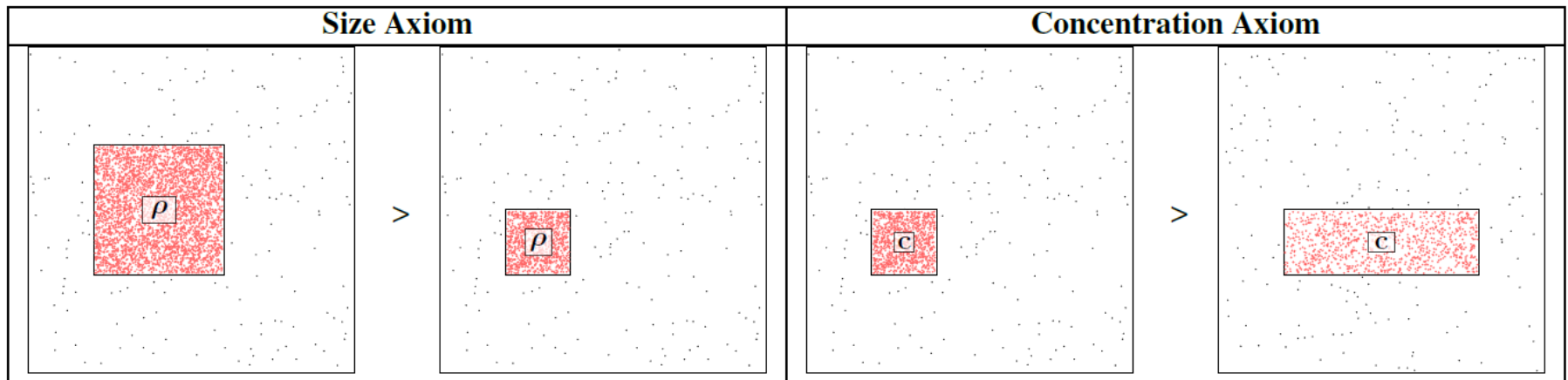
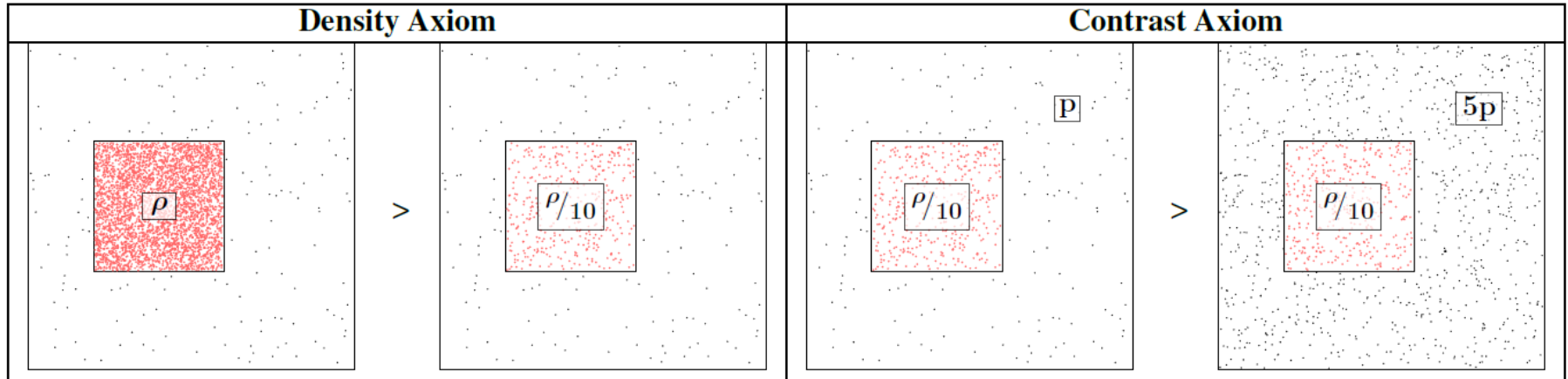
$$f\left(\begin{array}{c} n_1 \times n_2 \times n_3 \\ \text{mass } c \\ \text{density } \rho \end{array}\right)$$

VS

$$f\left(\begin{array}{c} n'_1 \times n'_2 \times n'_3 \\ \text{mass } c' \\ \text{density } \rho' \end{array}\right)$$

# Axioms: 1 to 4

$$c_1 > c_2 \iff f(\mathbf{n}, c_1, \mathbf{N}, C) > f(\mathbf{n}, c_2, \mathbf{N}, C)$$

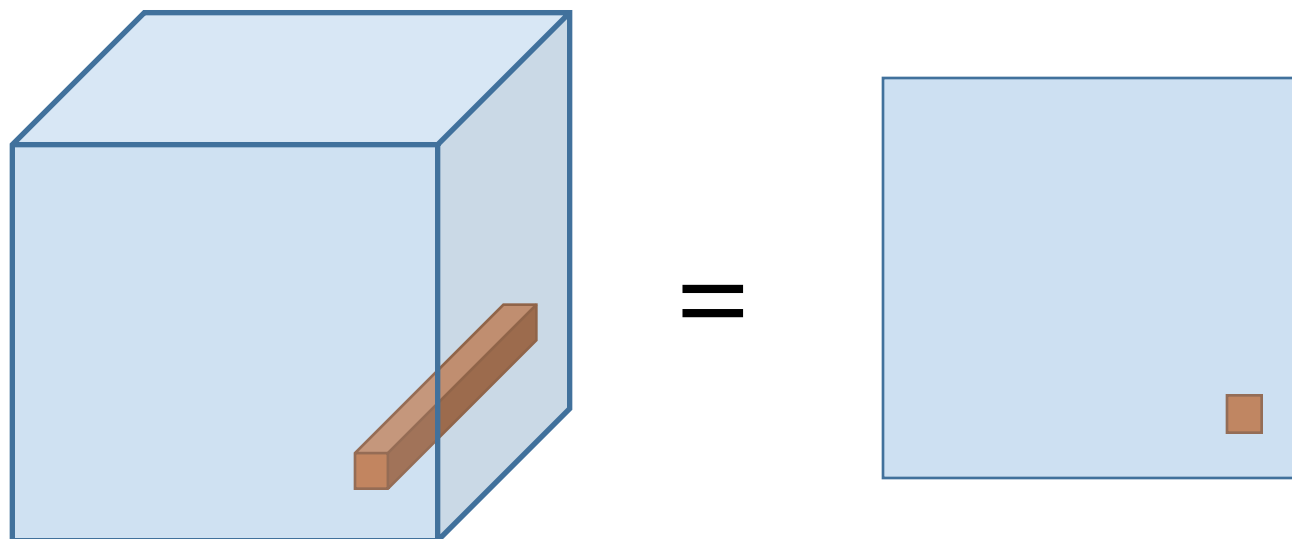


$$p_1 < p_2 \iff \hat{f}(\mathbf{n}, \rho, \mathbf{N}, p_1) > \hat{f}(\mathbf{n}, \rho, \mathbf{N}, p_2)$$

# Axiom 5: Cross Dimensions

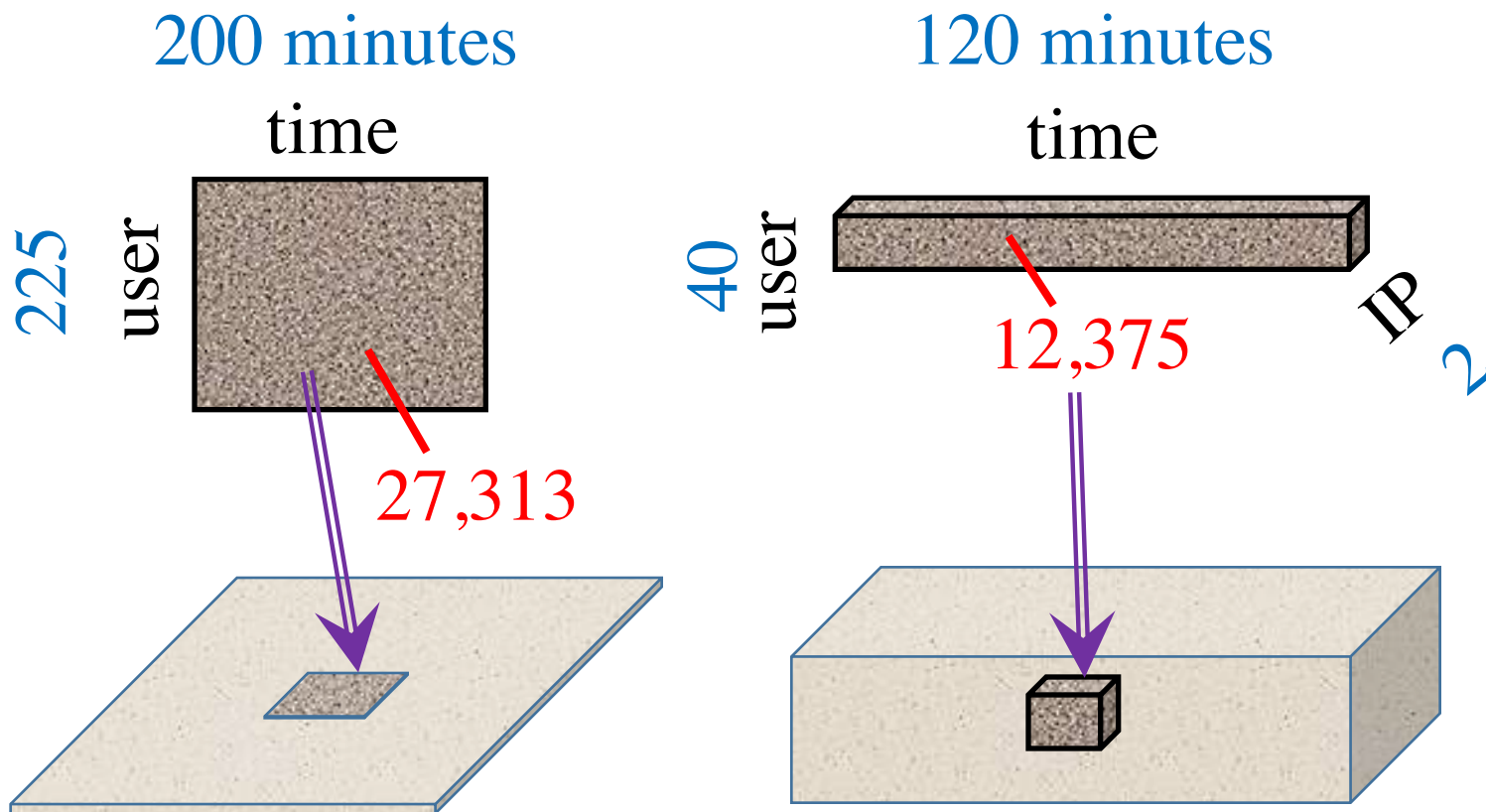
$$f_{K-1} \left( [n_k]_{k=1}^{K-1}, c, [N_k]_{k=1}^{K-1}, C \right) = f_K \left( ([n_k]_{k=1}^{K-1}, N_K), c, [N_k]_{k=1}^K, C \right)$$

Not including a mode is the same as including all values for that mode.



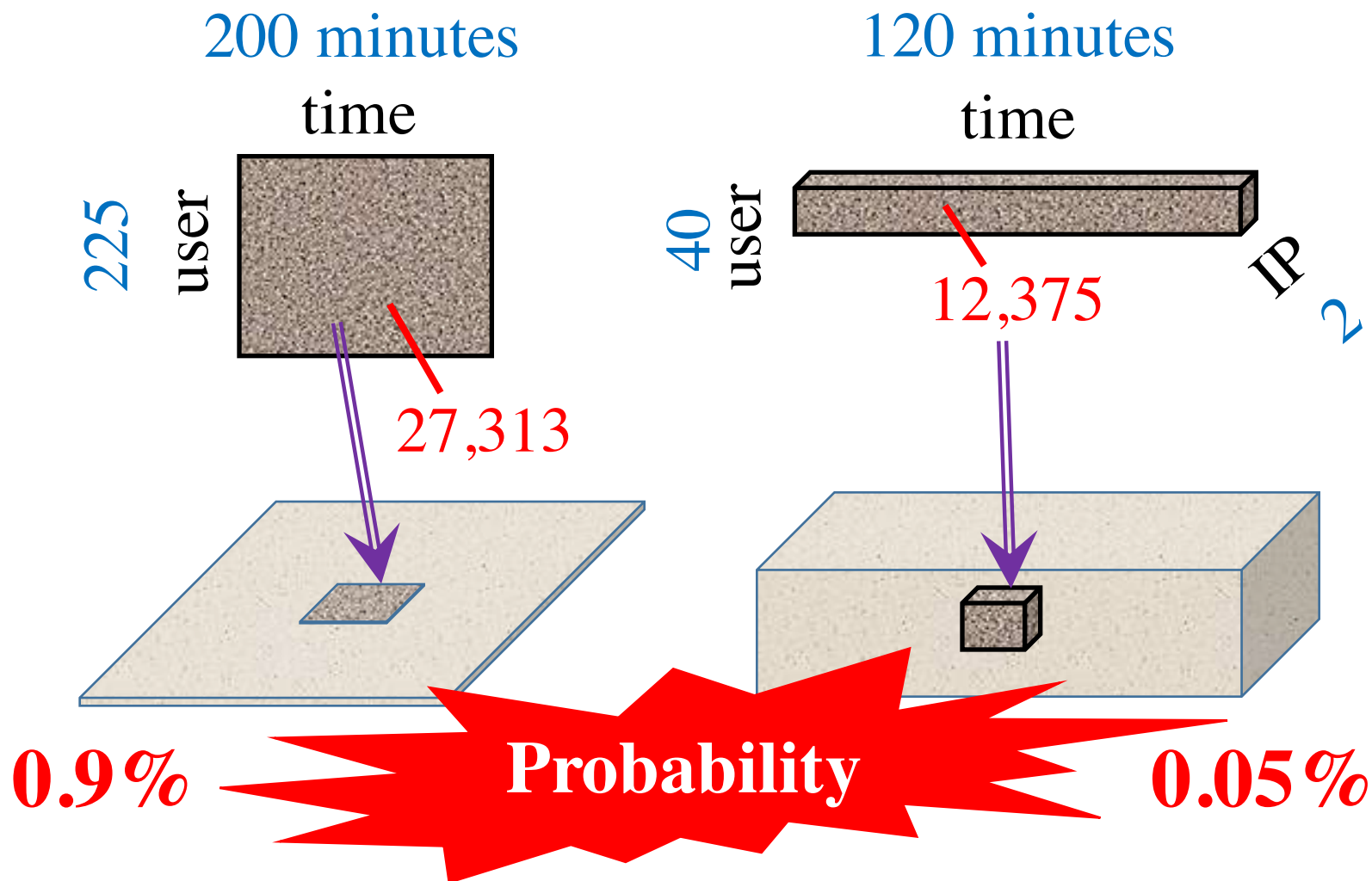
- ▶ New information (more modes) can only make our blocks more suspicious

# Scoring the Suspiciousness



**Q:** Which is more suspicious?

# Scoring the Suspiciousness





# A General Suspiciousness Metric

- Negative log likelihood of block's probability

$$f(n, c, N, C) = -\log [Pr(Y_n = c)]$$

**Lemma** Given an  $n_1 \times \dots \times n_K$  block of mass  $c$  in  $N_1 \times \dots \times N_K$  data of total mass  $C$ , the suspiciousness function is

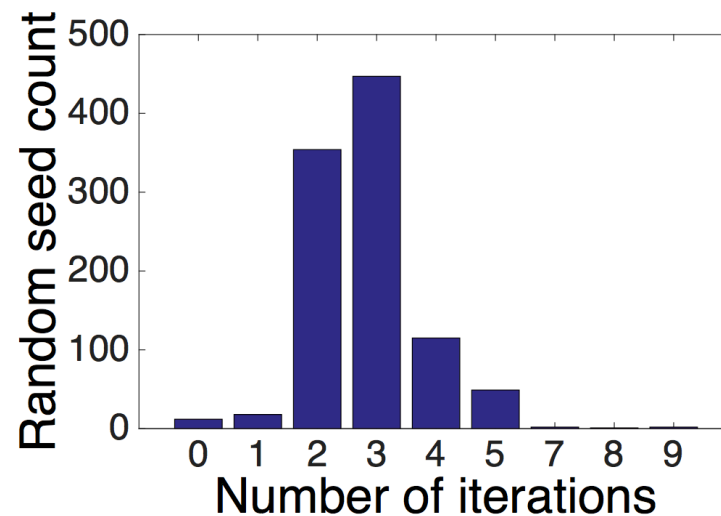
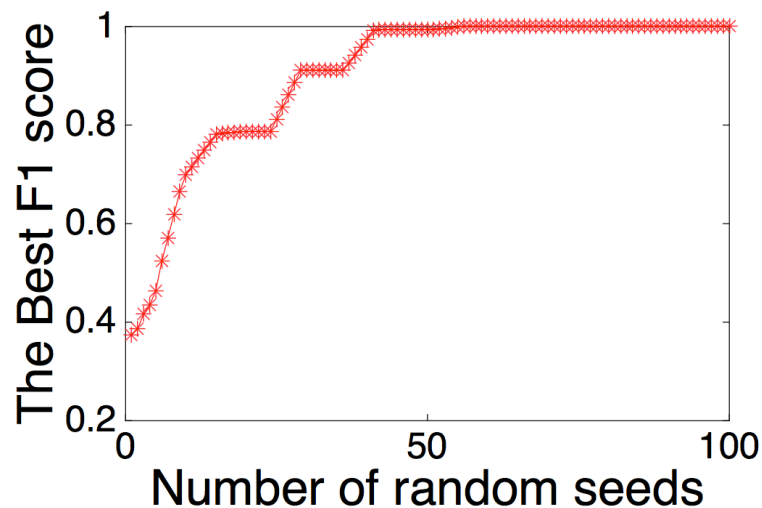
$$f(\mathbf{n}, c, \mathbf{N}, C) = c \left( \log \frac{c}{C} - 1 \right) + C \prod_{i=1}^K \frac{n_i}{N_i} - c \sum_{i=1}^K \log \frac{n_i}{N_i}$$

Using  $\rho$  as the block's density and  $p$  is the data's density, we have the simpler formulation

$$\hat{f}(\mathbf{n}, \rho, \mathbf{N}, p) = \left( \prod_{i=1}^K n_i \right) D_{KL}(\rho || p)$$

# CrossSpot Algorithm

- Local search to maximize the metric
  - Start with seed blocks
  - Parameter-free: iteratively update the blocks
  - Scalable: parallelize to multiple machines



# Advantages

		Axioms					
		Density	Size	Concentration	Contrast	Multi-modal	
Method		Scores Blocks	1	2	3	4	5
Metrics	<b>SUSPICIOUSNESS</b>	✓	✓	✓	✓	✓	✓
	Mass	✓	✓	✗	✗	✗	✓
	Density	✓	✓	✗	✓	✗	✗
	Average Degree [9]	✓	✓	✗	✗	✗	N/A
	Singular Value [10]	✓	✓	✓	✓	✗	✗
Methods	<b>CROSSPOT</b>	✓	✓	✓	✓	✓	✓
	Subgraph [30, 10, 36]	✓	✓	✓	✓	✗	N/A
	CopyCatch [6]	✓	✓	✓	✓	✗	N/A
	EigenSpokes [31]	✗	N/A				
	TrustRank [14, 8]	✗	N/A				
	BP [28, 1]	✗	N/A				

# Results: Dense Block Detection

## □ Synthetic data

□  $1,000 \times 1,000 \times 1,000$  of 10,000 random data

□ Block#1:  $30 \times 30 \times 30$  of 512 3 modes

□ Block#2:  $30 \times 30 \times 1,000$  of 512 2 modes

□ Block#3:  $30 \times 1,000 \times 30$  of 512 2 modes

□ Block#4:  $1,000 \times 30 \times 30$  of 512 2 modes

	Recall				Overall Evaluation		
	Block #1	Block #2	Block #3	Block #4	Precision	Recall	F1 score
HOSVD ( $r=20$ )	93.7%	29.5%	23.7%	21.3%	<b>0.983</b>	0.407	0.576
HOSVD ( $r=10$ )	91.3%	24.4%	18.5%	19.2%	0.972	0.317	0.478
HOSVD ( $r=5$ )	85.7%	10.0%	9.5%	11.4%	0.952	0.195	0.324
CROSSPOT	<b>100%</b>	<b>99.9%</b>	<b>94.9%</b>	<b>95.4%</b>	0.978	<b>0.967</b>	<b>0.972</b>

# Results: Tweeting Hashtags

User $\times$ hashtag $\times$ IP $\times$ minute	Mass $c$	Suspiciousness
$582 \times 3 \times 294 \times \mathbf{56,940}$	5,941,821	111,799,948
$188 \times 1 \times 313 \times \mathbf{56,943}$	2,344,614	47,013,868
$75 \times 1 \times 2 \times 2,061$	689,179	19,378,403

User ID	Time	IP address (city, province)	Tweet text with hashtag
USER-D	11-18 12:12:51	IP-1 (Deyang, Shandong)	<b>#Snow#</b> the Samsung GALAXY SII QQ Service customized version...
USER-E	11-18 12:12:53	IP-1 (Deyang, Shandong)	<b>#Snow#</b> the Samsung GALAXY SII QQ Service customized version...
USER-F	11-18 12:12:54	IP-2 (Zaozhuang, Shandong)	<b>#Snow#</b> the Samsung GALAXY SII QQ Service customized version...
USER-E	11-18 12:17:55	IP-1 (Deyang, Shandong)	<b>#Li Ning - a weapon with a hero#</b> good support activities!
USER-F	11-18 12:17:56	IP-2 (Zaozhuang, Shandong)	<b>#Li Ning - a weapon with a hero#</b> good support activities!
USER-D	11-18 12:18:40	IP-1 (Deyang, Shandong)	<b>#Toshiba Bright Daren#</b> color personality test to find out your sense...
USER-E	11-18 17:00:31	IP-2 (Zaozhuang, Shandong)	<b>#Snow#</b> the Samsung GALAXY SII QQ Service customized version...
USER-D	11-18 17:00:49	IP-2 (Zaozhuang, Shandong)	<b>#Toshiba Bright Daren#</b> color personality test to find out your sense...
USER-F	11-18 17:00:56	IP-2 (Zaozhuang, Shandong)	<b>#Li Ning - a weapon with a hero#</b> good support activities!

# Results: Network Attacks

	#	Src-IP $\times$ dst-IP $\times$ port $\times$ second	Mass $c$	Suspiciousness
CROSSPOT	1	$411 \times 9 \times 6 \times \mathbf{3,610}$	47,449	552,465
	2	$533 \times 6 \times 1 \times \mathbf{3,610}$	30,476	400,391
	3	$5 \times 5 \times 2 \times \mathbf{3,610}$	18,881	317,529
	4	$11 \times 7 \times 7 \times \mathbf{3,610}$	20,382	295,869
HOSVD	1	$15 \times 1 \times 1 \times 1,336$	4,579	80,585
	2	$1 \times 2 \times 2 \times 1,035$	1,035	18,308
	3	$1 \times 1 \times 1 \times 1,825$	1,825	34,812
	4	$1 \times 13 \times 6 \times 181$	1,722	29,224



# Summary

- ❑ Ill-gotten Facebook Likes, Zombie Followers
- ❑ **Observations, Representations, Models**
  - ❑ **CopyCatch:** Catching ill-gotten Likes by core search
  - ❑ **LockInfer:** Adding seed selection before search
  - ❑ **CatchSync:** Catching smart zombie followers with high recall (recovering power-law distributions)
  - ❑ **CrossSpot:** Defining suspiciousness across dimensions



# Acknowledgement





# References

- D. Blei, A. Ng, and M. Jordan. “Latent dirichlet allocation.” JMLR, 2003.
- J. Herlocker, J. Konstan, L. Terveen, J. Riedl. “Evaluating collaborative filtering recommender systems.” ACM TOIS, 2004.
- Y. Koren, R. Bell, C. Volinsky. “Matrix factorization techniques for recommender systems.” Computer, 2009.
- Y. Koren. “Factorization meets the neighborhood: A multifaceted collaborative filtering model.” KDD, 2008.
- Y. Koren. “Collaborative filtering with temporal dynamics.” CACM, 2010.
- M. Balabanovic and Y. Shoham. “FAB: Content-based, collaborative recommendation.” CACM, 1997.
- N. Liu and Q. Yang. “Eigenrank: A ranking-oriented approach to collaborative filtering.” SIGIR, 2008.
- N. Liu, M. Zhao, and Q. Yang. “Probabilistic latent preference analysis for collaborative filtering.” CIKM, 2009.

# References

H. Ma, H. Yang, M. Lyu, and I. King. “Sorec: Social recommendation using probabilistic matrix factorization.” CIKM, 2008.

H. Ma, T. Zhou, M. Lyu, and I. King. “Improving recommender systems by incorporating social contextual information.” ACM TOIS, 2011.

H. Ma, D. Zhou, C. Liu, M. Lyu, and I. King. “Recommender systems with social regularization.” WSDM, 2011.

J. Leskovec, A. Singh, and J. Kleinberg. “Patterns of influence in a recommendation network.” PAKDD, 2006.

P. Massa and A. Paolo. “Trust-aware recommender systems.” RecSys, 2007.

M. Jamali and E. Martin. “TrustWalker: A random walk model for combining trust-based and item-based recommendation.” KDD, 2009.

H. Ma, I. King, and M. Lyu. “Learning to recommend with social trust ensemble.” SIGIR, 2009.

H. Ma, I. King, and M. Lyu. “Learning to recommend with explicit and implicit social relations.” ACM TIST, 2011.



# References

- M. Faloutsos, P. Faloutsos, and C. Faloutsos. “On power-law relationships of the internet topology.” SIGCOMM, 1999.
- A. Broder, R. Kumar, F. Maghoul, P. Raghavan, S. Rajagopalan, R. Stata, A. Tomkins, and J. Weiner. “Graph structure in the web.” Computer Networks, 2000.
- F. Chung and L. Lu. “The average distances in random graphs with given expected degrees.” PNAS, 2002.
- J. Kleinberg. “Authoritative sources in a hyperlinked environment.” JACM, 1999.
- H. Kwak, C. Lee, H. Park, and S. Moon. “What is Twitter, a social network or a news media?” WWW, 2010.
- B. Hooi, H.A. Song, A. Beutel, N. Shah, K. Shin, and C. Faloutsos. “Fraudar: Bounding graph fraud in the face of camouflage.” KDD, 2016.
- C. Aggarwal and J. Han. “Frequent pattern mining.” Springer, 2014.
- J. Han, J. Pei, B. Mortazavi-Asl, Q. Chen, U. Dayal, and M. Hsu. “FreeSpan: Frequent pattern-projected sequential pattern mining.” KDD, 2000.

# References

- X. Yan and J. Han. “gspan: Graph-based substructure pattern mining.” ICDM, 2003.
- X. Yan and J. Han. “CloseGraph: Mining closed frequent graph patterns.” KDD, 2003.
- Y. Sun, J. Han, X. Yan, P.S. Yu, and T. Wu. “PathSim: Meta path-based top-k similarity search in heterogeneous information networks.” VLDB, 2011.
- Y. Sun, Y. Yu, and J. Han. “Ranking-based clustering of heterogeneous information networks with star network schema.” KDD, 2009.
- Y. Sun, J. Han, P. Zhao, Z. Yin, H. Cheng, and T. Wu. “RankClus: Integrating clustering with ranking for heterogeneous information network analysis.” EDBT, 2009.
- Y. Sun, R. Barber, M. Gupta, C. Aggarwar, and J. Han. “Co-author relationship prediction in heterogeneous bibliographic networks.” ASONAM, 2011.
- A. El-Kishky, Y. Song, C. Wang, C.R. Voss, and J. Han. “Scalable topical phrase mining from text corpora.” VLDB, 2014.
- J. Liu, J. Shang, C. Wang, X. Ren, and J. Han. “Mining quality phrases from massive text corpora.” SIGMOD, 2015.

# References

X. Ren, A. El-Kishky, C. Wang, F. Tao, C.R. Voss, and J. Han. “Effective entity recognition and typing by relation phrase-based clustering.” KDD, 2015.

X. Ren, W. He, M. Qu, C.R. Voss, H. Ji, and J. Han. “Label noise reduction in entity typing by heterogeneous partial-label embedding.” KDD, 2016.

C. Wang, M. Danilevsky, N. Desai, Y. Zhang, P. Nguyen, T. Taula, and J. Han. “A phrase mining framework for recursive construction of a topical hierarchy.” KDD, 2013.

E.E. Papalexakis, C. Faloutsos, N.D. Sidiropoulos. “ParCube: Sparse parallelizable tensor decompositions.” PKDD, 2012.

D. Koutra, U. Kang, J. Vreeken, and C. Faloutsos. “VOG: Summarizing and understanding large graphs.” SDM, 2014.

R. Gupta, A. Halevy, X. Wang, S.E. Whang, and F. Wu. “Biperpedia: An ontology for search applications.” VLDB, 2014.

M. Yahya, S. Whang, R. Gupta, and A. Halevy. “ReNoun: Fact extraction for nominal attributes.” EMNLP, 2014.

A. Halevy, N. Noy, S. Sarawagi, S.E. Whang, and X. Yu. “Discovering structure in the universe of attribute names.” WWW, 2016.

# References

Q. Li, Y. Li, J. Gao, B. Zhao, W. Fan, and J. Han. “Resolving conflicts in heterogeneous data by truth discovery and source reliability estimation.” SIGMOD, 2014.

Q. Li, Y. Li, J. Gao, L. Su, B. Zhao, M. Demirbas, W. Fan, and J. Han. “A confidence-aware approach for truth discovery on long-tail data.” VLDB, 2014.

F. Ma, Y. Li, Q. Li, M. Qiu, J. Gao, S. Zhi, L. Su, B. Zhao, H. Ji, and J. Han. “Faitcrowd: Fine grained truth discovery for crowdsourced data aggregation.” KDD, 2015.

Y. Li, J. Gao, C. Meng, Q. Li, L. Su, B. Zhao, W. Fan, and J. Han. “A survey on truth discovery.” KDD Explorations Newsletter, 2016.

S. Zhi, B. Zhao, W. Tong, J. Gao, D. Yu, H. Ji, and J. Han. “Modeling truth existence in truth discovery.” KDD, 2015.

S. Kumar, R. West, and J. Leskovec. “Disinformation on the Web: Impact, characteristics, and detection of Wikipedia hoaxes.” WWW, 2016.

S. Kumar, F. Spezzano, and V.S. Subrahmanian. “Identifying malicious actors on social media.” ASONAM, 2016. (tutorial)



# Thank you!

**Data-Driven Behavioral Analytics:  
Observations, Representations and Models**